

**PCT**WORLD INTELLECTUAL PROPERTY ORGANIZATION  
International Bureau

## INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(51) International Patent Classification <sup>6</sup> :**H04L 12/28, 12/56****A1**

(11) International Publication Number:

**WO 99/00939**

(43) International Publication Date:

**7 January 1999 (07.01.99)**(21) International Application Number: **PCT/US98/13365**(22) International Filing Date: **25 June 1998 (25.06.98)**

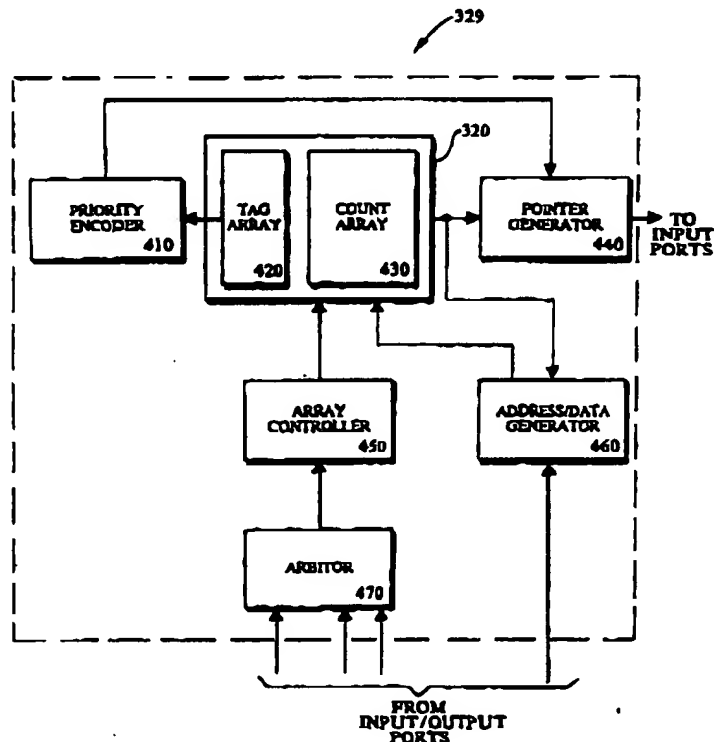
(30) Priority Data:

08/885,118	30 June 1997 (30.06.97)	US
08/987,914	9 December 1997 (09.12.97)	US

(71) Applicant: **SUN MICROSYSTEMS, INC. [US/US]; 901 San Antonio Road, Palo Alto, CA 94303 (US).**(72) Inventors: **MULLER Shimon; Apartment D, 983 La Mesa Terrace, Sunnyvale, CA 94086 (US). HENDEL, Ariel; 7537 Newcastle Drive, Cupertino, CA 95014 (US). TANGIRALA, Ravi; 2065 California Avenue, Mountain View, CA 94040 (US). BERG, Curt; 1370 Ensenada Way, Los Altos, CA 94024 (US).**(74) Agents: **HYMAN, Eric, S. et al.; Blakely, Sokoloff, Taylor & Zafman, 7th floor, 12400 Wilshire Boulevard, Los Angeles, CA 90025-1026 (US).**(81) Designated States: **JP, European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE).****Published***With International search report.*(54) Title: **SHARED MEMORY MANAGEMENT IN A SWITCHED NETWORK ELEMENT**

(57) Abstract

A method and apparatus for shared memory management (220) in a switched network element (100) are provided. According to one aspect of the present invention, a shared memory manager (220) for a packet forwarding device (140) includes a pointer memory (320) having stored therein information regarding buffer usage (e.g., usage counts) for each of a number of buffers in a shared memory (230). An encoder (410) is coupled to the pointer memory (320) for generating an output which indicates a set of buffers that contains a free buffer. The shared memory manager (220) further includes a pointer generator (440) that is coupled to the encoder (410) for locating a free buffer in the set of buffers. The pointer generator (440) is further configured to produce a pointer to the free buffer based upon the output of the encoder (410) and the free buffer's location within the set of buffers. According to another aspect of the present invention, a packet forwarding device (100) includes a number of output ports (117) for transmitting packets onto a network and a number of input ports (117) coupled to the output ports for receiving packets from the network, buffering the packets, and forwarding the packets to one or more of the output ports (117). The packet forwarding device (100) also includes a shared memory (230) that is segmented into buffers (140) for temporarily buffering the packets. No more than one copy of a given packet is ever stored in the shared memory (230). The packet forwarding device (100) further includes a shared memory manager (220) which dynamically allocates buffers on behalf of the input ports (117) and tracks ownership counts for each of the buffers.



WO-99/00939

PCT/US98/13365

- 1 -

## SHARED MEMORY MANAGEMENT IN A SWITCHED NETWORK ELEMENT

5 CROSS-REFERENCES TO RELATED APPLICATIONS

The present application is a continuation-in-part of co-pending U.S. Patent Application Number 08/885,118, entitled, "Shared Memory Management in a Switched Network Element" filed on June 30, 1997, attorney docket number 082225.P2354.

10 FIELD OF THE INVENTION

The invention relates generally to the field of packet forwarding in computer networking devices. More particularly, the invention relates to shared memory management in a switched network element.

15 BACKGROUND OF THE INVENTION

An increasing number of users are requiring increased bandwidth from existing networks due to multimedia applications for accessing the Internet and World Wide Web, for example. Therefore, future networks must be able to support a very high bandwidth and a large number of users. Furthermore, such networks should be able to support  
20 multiple traffic types such as data, voice, and video which typically require different bandwidths.

Statistical studies indicate that the network domain, i.e., a group of interconnected local area networks (LANs), as well as the number of individual end-stations connected to each LAN, will grow at ever increasing rates in the future. Thus, more network  
25 bandwidth and more efficient use of resources is needed to meet these increasing rates.

WO 99/00939

PCT/US98/13365

- 2 -

A common source of inefficiency in prior switched network elements is the memory management mechanism for packet buffering. Packet buffering is typically required in a switched network element to avoid packet loss. One potential cause of congestion is a speed mismatch between an input port and an output port. For example, if

5 a fast input port (e.g., 1,000Mb/s) forwards traffic to a slow output port (e.g., 10Mb/s), the slower output port will not be able to transmit packets onto the network as fast as it is receiving packets from the faster input port. Thus, packets must be buffered or they will be dropped. Particular traffic patterns may also result in congestion. The traffic patterns crossing the switched network element may be such that several input ports need to

10 forward data to the same output port, for example. As a result, temporary congestion on that output port may occur. Further, multicast traffic arriving at one or more input ports may need to be forwarded to many output ports. This causes traffic multiplication which may also result in temporary congestion on one or more output ports. Finally, competition for common resources may contribute to congestion. For example, common resources

15 required for packet forwarding may cause incoming traffic to accumulate on one or more input ports. Packets may need to be buffered at a particular input port while another input port is accessing a particular common resource such as a forwarding database.

Typically, one of two approaches is employed to achieve the required packet buffering. The first approach, input port buffering, associates packet (buffer) memory to

20 input ports for temporarily storing packet data until it can be forwarded to the appropriate output port(s). The second approach, output port buffering, associates packet memory to the output port for temporary storage of packet data until it can be transmitted onto the attached link.

WO 99/00939

PCT/US98/13365

- 3 -

A major architectural challenge in implementing a high performance switched network element is the provision of just the right amount of packet buffering for each port. An inadequate amount of packet memory, even on only one of the ports, may have serious performance implications for the entire switch. On the other hand, too much buffering will unnecessarily increase the cost of the switching fabric with no incremental benefit. Due to the difficulty of estimating buffering requirements for each port, many implementations either cost too much or do not perform very well, or both.

Based on the foregoing, it should be apparent that one candidate for improved efficiency is the memory management mechanism of a networking device. Further, recognizing the intrinsic efficiency of sharing resources and the bursty nature of network traffic, it is desirable to utilize a dynamic packet memory management scheme to facilitate sharing of a common packet memory among all input/output ports for packet buffering.

WO 99/00939

PCT/US98/13365

- 4 -

SUMMARY OF THE INVENTION

A method and apparatus for shared memory management in a switched network element is described. According to one aspect of the present invention, a shared memory manager for a packet forwarding device includes a pointer memory having stored therein  
5 information regarding buffer usage for each of a number of buffers in a shared memory. An encoder is coupled to the pointer memory. The encoder is configured to generate an output which indicates a set of buffers that contains one or more free buffers. The shared memory manager further includes a pointer generator. The pointer generator is coupled to the encoder and is configured to locate a free buffer in the set of buffers. The pointer  
10 generator is further configured to produce a pointer to the free buffer based upon the output of the encoder and the free buffer's location within the set of buffers.

According to another aspect of the present invention, a packet forwarding device includes a number of output ports for transmitting packets onto a network and a number of input ports coupled to the output ports for receiving packets from the network, buffering  
15 the packets, and forwarding the packets to one or more of the output ports. The packet forwarding device also includes a shared memory coupled to the output ports and the input ports. The shared memory is segmented into a number of buffers for temporarily buffering the packets. However, at any given time, no more than one copy of a given packet is stored in the shared memory. The packet forwarding device further includes a  
20 shared memory manager coupled to the input ports and to the output ports. The shared memory manager dynamically allocates buffers on behalf of the input ports and tracks ownership counts for each of the buffers based upon information provided by the input ports and the output ports.

WO 99/00939

PCT/US98/13365

- 5 -

According to yet another aspect of the present invention, a method is provided for packet forwarding. The method includes dynamically allocating one or more buffer pointers that identify one or more buffers in a shared memory. When a packet is received the packet is stored in the one or more buffers. Then, the buffer pointers are transferred  
5 based upon a forwarding decision. Finally, the packet is transmitted after retrieving the packet from the buffers.

Other features of the present invention will be apparent from the accompanying drawings and from the detailed description which follows.

WO 99/00939

PCT/US98/13365

- 6 -

**BRIEF DESCRIPTION OF THE DRAWINGS**

The present invention is illustrated by way of example, and not by way of limitation, in the figures of the accompanying drawings and in which like reference numerals refer to similar elements and in which:

5

Figure 1 illustrates a switch according to one embodiment of the present invention.

Figure 2 is a simplified block diagram of an exemplary switching element that may be utilized in the switch of Figure 1.

10 Figure 3A is a logical view of the shared memory of Figure 2 according to one embodiment of the present invention.

Figure 3B is a block diagram of the shared memory manager of Figure 2 according to one embodiment of the present invention.

Figure 4 is a block diagram of the buffer tracking process of Figure 3B according to one embodiment of the present invention.

15 Figure 5 is a flow diagram illustrating buffer allocation processing according to one embodiment of the present invention.

Figure 6 is a flow diagram illustrating buffer ownership transfer processing according to one embodiment of the present invention.

20 Figure 7 is a flow diagram illustrating buffer return processing according to one embodiment of the present invention.

WO 99/00939

PCT/US98/13365

- 7 -

DETAILED DESCRIPTION

A method and apparatus for shared memory management in a switched network element is described. In the following description, for the purposes of explanation, numerous specific details are set forth in order to provide a thorough understanding of the present invention. It will be apparent, however, to one skilled in the art that the present invention may be practiced without some of these specific details. In other instances, well-known structures and devices are shown in block diagram form.

The present invention includes various steps, which will be described below. While the steps of the present invention are preferably performed by the hardware components described below, the steps may alternatively be embodied in machine-executable instructions stored in a machine-readable medium, such as a memory, CD-ROM, diskette or other storage medium, which may be used to cause a general-purpose or special-purpose processor programmed with the instructions to perform the steps. Further, embodiments of the present invention will be described with reference to a high speed Ethernet switch. However, the method and apparatus described herein are equally applicable to other types of network devices and protocols.



WO 99/00939

PCT/US98/13365

- 8 -

## AN EXEMPLARY NETWORK ELEMENT

An overview of one embodiment of a network element that operates in accordance with the teachings of the present invention is illustrated in Figure 1. The network element  
5 is used to interconnect a number of nodes and end-stations in a variety of different ways. In particular, an application of a multi-layer distributed network element (MLDNE) would be to forward packets according to predefined protocols over a homogenous data link layer such as the IEEE 802.3 standard, also known as Ethernet. Other protocols can also be used.

10 The MLDNE's distributed architecture can be configured to forward message traffic in accordance with a number of known or future forwarding algorithms. In a preferred embodiment, the MLDNE is configured to handle message traffic using the Internet suite of protocols, and more specifically the Transmission Control Protocol (TCP) and the Internet Protocol (IP) over the Ethernet LAN standard and medium access control  
15 (MAC) data link layer. TCP is also referred to here as a Layer 4 protocol, while the IP is referred to repeatedly as a Layer 3 protocol. For purposes of discussion, references to Layers herein typically refer to the Open Systems Interconnection (OSI) seven layer model created by the International Organization for Standardization (ISO).

In one embodiment of the MLDNE, a network element is configured to implement  
20 packet forwarding functions in a distributed manner, i.e., different parts of a function are performed by different subsystems in the MLDNE, while the final result of the functions remains transparent to the nodes and end-stations. As will be appreciated from the discussion below and the diagram in Figure 1, the MLDNE has a scalable architecture which allows the designer to predictably increase the number of external connections by

WO 99/00939

PCT/US98/13365

- 9 -

adding additional subsystems, thereby allowing greater flexibility in defining the MLDNE as a stand alone router.

As illustrated in block diagram form in Figure 1, the MLDNE 101 contains a number of subsystems 110 that are interconnected using a number of internal links 141 to create a larger switch. According to one embodiment, the subsystems 110 may be fully meshed by providing at least one internal link between any two subsystems. Each subsystem 110 includes a switching element 100 coupled to a forwarding and filtering memory 140, also referred to as a forwarding database. The forwarding and filtering database may include a forwarding memory 113 and an associated memory 114. The forwarding memory (or database) 113 stores an address table used for matching with the headers of received packets. The associated memory (or database) stores data associated with each entry in the forwarding memory that is used to identify forwarding attributes for forwarding the packets through the MLDNE. A number of external ports (not shown) having input and output capability interface the external connections 117. In one embodiment, each subsystem supports multiple Gigabit Ethernet ports (The term Gigabit Ethernet, as used herein shall apply to networks employing Carrier Sense, Multiple Access with Collision Detection (CSMA/CD) as the medium access method, generally operating at a signaling rate of 1,000 Mb/s over various media types and transmitting Ethernet formatted or Institute of Electrical and Electronic Engineers (IEEE) standard 802.3 formatted data packets), Fast Ethernet ports (The term Fast Ethernet, as used herein shall apply to networks employing CSMA/CD as the medium access method, generally operating at a signaling rate of 100 Mb/s over various media types and transmitting Ethernet formatted or IEEE standard 802.3 formatted data packets) and Ethernet ports (The term Ethernet, as used herein shall apply to networks employing CSMA/CD as the

WO 99/00939

PCT/US98/13365

- 10 -

medium access method, generally operating at a signaling rate of 10 Mb/s over various media types and transmitting Ethernet formatted or IEEE standard 802.3 formatted data packets). Internal links 141 are used to couple internal ports (not shown). Using the internal links, the MLDNE can connect multiple switching elements together to form a multigigabit switch.

The MLDNE 101 further includes a central processing system (CPS) 160 that is coupled to the individual subsystem 110 through a communication bus 151, such as the peripheral components interconnect (PCI). PCI is mentioned merely as an exemplary communication bus, those of ordinary skill in the art will appreciate the type of bus may vary for different implementations. The CPS 160 includes a central processing unit (CPU) 161 coupled to a central memory 163. Central memory 163 includes a copy of the data contained in the individual forwarding memories 113 of the various subsystems 110. The CPS 160 has a direct control and communication interface to each subsystem 110 and provides some centralized communication and control between switching elements 100.

#### AN EXEMPLARY SWITCHING ELEMENT

Figure 2 is a simplified block diagram illustrating an exemplary architecture of the switching element of Figure 1. The switching element 100 depicted includes a central processing unit (CPU) interface 215, a switch fabric block 210, a network interface 205, a cascading interface 225, and a shared memory manager 220.

Packets may enter or leave the network switching element 100 through any one of the three interfaces 205, 215, or 225. In brief, the network interface 205 operates in accordance with a network communication protocol, such as Ethernet, to receive packets

- 11 -

from a network (not shown) and to transmit packets onto the network via one or more input ports and output ports 206, respectively. An optional cascading interface 225 may include one or more internal links 226 for interconnecting switching elements 100 to create larger switches. For example, each switching element 100 may be connected together  
5 with other switching elements 100 in a full mesh topology to form a multi-layer switch as described above. Alternatively, a switch may comprise a single switching element 100 with or without the cascading interface 225.

The CPU 161 may give commands or packets to the network switching element 100 via the CPU interface 215. In this manner, one or more software processes running  
10 on the CPU 161 may manage the entries in the external forwarding and filtering database 140, such as adding new entries and invalidating unwanted entries. In alternative embodiments, however, the CPU 161 may be provided with direct access to the forwarding and filtering database 140. In any event, for purposes of packet forwarding, the CPU port of the CPU interface 215 resembles a generic port into the switching element  
15 100 and may be treated as if it were simply another external network interface port. However, since access to the CPU port occurs over a bus such as a peripheral components interconnect (PCI) bus, the CPU port does not need any media access control (MAC) functionality.

Returning to the network interface 205, the two main tasks of input packet  
20 processing and output packet processing will now briefly be described. Input packet processing may be performed by one or more input ports of the network interface 205. Input packet processing includes the following: (1) receiving and verifying incoming Ethernet packets, (2) modifying packet headers when appropriate, (3) requesting buffer pointers from the shared memory manager 220 for storage of incoming packets, (4)

WO 99/00939

PCT/US98/13365

- 12 -

requesting forwarding decisions from the switch fabric block 210, (5) transferring the incoming packet data to the shared memory manager 220 for temporary storage in an external shared memory 230, and (5) upon receipt of a forwarding decision, forwarding the buffer pointer(s) to the output port(s) 206 indicated by the forwarding decision.

- 5 Output packet processing may be performed by one or more output ports 206 of the network interface 205. Output processing may include requesting packet data from the shared memory manager 220, transmitting packets onto the network, and requesting deallocation of buffer(s) after packets have been transmitted.

- 10 The network interface 205, the CPU interface 215, and the cascading interface 225 are coupled to the shared memory manager 220 and the switch fabric block 210. The shared memory manager 220 provides an efficient centralized interface to the external shared memory 230 for buffering of incoming packets. The switch fabric block 210 includes a search engine and learning logic for searching and maintaining the forwarding and filtering database 140 with the assistance of the CPU 161.

- 15 The switch fabric block 210 includes a search engine that provides access to the forwarding and filtering database 140 on behalf of the interfaces 205, 215, and 225. Packet header matching, learning, packet forwarding, filtering, and aging are exemplary functions that may be performed by the switch fabric block 210. Each input port 206 is coupled with the switch fabric block 210 to receive forwarding decisions for received  
20 packets. The forwarding decision indicates the outbound port(s) (e.g., external network port or internal cascading port) upon which the corresponding packet should be transmitted. Additional information may also be included in the forwarding decision to support hardware routing such as a new MAC destination address (DA) for MAC DA

WO 99/00939

PCT/US98/13365

- 13 -

replacement. Further, a priority indication may also be included in the forwarding decision to facilitate prioritization of packet traffic through the switching element 100.

In the present embodiment, Ethernet packets are centrally buffered by the shared memory manager 220. The shared memory manager 220 interfaces every input port and output port 206 and performs dynamic memory allocation and deallocation on their behalf, respectively. During input packet processing, one or more buffers are allocated in the external shared memory 230 and an incoming packet is stored by the shared memory manager 220 responsive to commands received from the network interface 205, for example. Subsequently, during output packet processing, the shared memory manager 220 retrieves the packet from the external shared memory 230 and deallocates buffers that are no longer in use. Because multiple ports can own a given buffer, to assure no buffers are released until all output ports 206 have completed transmission of the data stored therein, the shared memory manager 220 preferably also tracks buffer ownership.

#### 15 PACKET SWITCHING OVERVIEW

According to one embodiment of the present invention, the switching element 100 of the present invention provides wire speed routing and forwarding of Ethernet, Fast Ethernet, and Gigabit Ethernet packets among the three interfaces 215, 205, and 225. By "wire speed" what is meant is the forwarding decision for a packet received on a given input port 206 is complete before the next packet arrives at that input port 206.

Forwarding is performed by passing pointers from input ports to output ports 206. The shared memory manager 220 provides a level of indirection which is exploited by the input and output ports 206 by locally storing pointers to buffers that contain packet data rather than locally storing the packet data itself. For example, input and output queues

WO 99/00939

PCT/US98/13365

- 14 -

may be maintained at input and output ports 206 respectively for temporarily storing pointers during input and output packet processing. The memory for buffering incoming packets is allocated from a common pool of memory (e.g., the shared memory 230) that is shared by all the input ports and output ports 206 of the switching element 100.

5 Briefly, the packet forwarding process begins with a packet being received at one of the switching element's input ports 206. It is important to note that, by keeping a predetermined number of buffer pointers on hand to allow immediate storage of received packet data, input ports 206 are always prepared to receive the next packet. These buffer pointers may be preallocated during the switching element's 100 initialization and  
10 subsequently requested from the shared memory manager 220 when the number of pointers falls below a predetermined threshold. Returning to the present example, the a portion of the received packet may be buffered temporarily at the input port 206 while a determination is made regarding the output port(s) 206 to which the packet is to be forwarded. Packets that are to be filtered, therefore, need not be stored in the shared  
15 memory 230.

After a forwarding decision is received for a particular packet, the input port 206 transfers ownership of the one or more buffers corresponding to the packet to the appropriate output port(s) 206. The transfer of ownership includes the input port 206 notifying the shared memory manager 220 of the number of output ports 206 that should  
20 transmit the packet and the input port 206 forwarding the appropriate pointers to those output ports 206.

Upon receipt of a buffer pointer, an output port 206 stores the pointer in an output queue until it can be transmitted onto the attached link. When the output port 206 is finished transmitting packet data from a particular buffer it notifies the shared memory

- 15 -

manager 220 that it is finished with the buffer. The shared memory manager 220 then updates its internal counts used for tracking the number of buffer owners and returns the buffer to the free pool if appropriate (e.g., the buffer is no longer in any output queues).

From the above overview, it should be appreciated that the use of buffer pointers  
5 reduces forwarding to the transfer of one or more buffer pointers from an input port 206 to one or more output ports 206. Further, flooding and the processing of multicast packets are made more efficient because packet data need not be duplicated. In fact, regardless of the number of output ports on which a particular packet is to be forwarded, only a single copy of the packet data will ever exist in the shared memory 230. Thus, as one advantage  
10 of the present embodiment, the architecture gracefully scales by accommodating an increased number of ports without requiring a proportionate increase in buffer memory.

#### SHARED MEMORY ORGANIZATION

Prior switching elements may have a fixed amount of memory associated with each  
15 port, resulting in inefficient memory allocation and buffering that is not related to the actual amount of traffic through a given port. Further, since the buffer memory is distributed, the logic for buffer management is duplicated for each port. In contrast, the shared memory manager 220 provides an efficient centralized interface to a shared pool of packet memory for buffering of incoming packets. Moreover, the memory management  
20 mechanism provided by the present invention is designed to achieve efficient allocation of per port buffering that is proportional to the amount of traffic through a given port. According to one embodiment, this proportional buffering is achieved by employing shared memory 230 in combination with a dynamic buffer allocation scheme. The shared memory 230 is a pool of buffers that is used for temporary storage of packet data en route



WO 99/00939

PCT/US98/13365

- 16 -

from an inbound interface (e.g., an input port 206 in the network interface 205, the cascading interface 225, or the CPU interface 215) to one or more outbound interfaces (e.g., an output port 206 in the network interface 205, the cascading interface 225, or the CPU interface 215). Essentially, the shared memory 230 serves as an elasticity buffer for  
5 adapting between the incoming and outgoing bandwidth requirements.

At this point, it may be useful to discuss the tradeoffs among certain shared memory parameters including buffer size, address space, and output/input pointer queue sizes. For example, a larger buffer size will more likely accommodate a full packet rather than a portion of a packet. However, potentially more buffer memory will be consumed  
10 when the packet size is not an integer multiple of the buffer size. Smaller buffer sizes, on the other hand, conserve memory in this situation due to the finer resolution. However, more addresses may be required to uniquely identify the buffers and each packet potentially requires more buffers for storage. Additionally, more pointers may need to be  
15 queued at both the input and output ports 206 as a result of increasing the number of buffers per packet. Further, if the environment is not known in advance, it is desirable to provide programmable resources, thereby allowing buffer sizes, the shared memory size, queue sizes, and other parameters to be optimized for a particular implementation. For example, in an Ethernet implementation a buffer size of 512 bytes will typically result in the use of one to three buffers per packet.

20 According to one embodiment of the present invention, the shared memory manager 220 includes a buffered architecture that utilizes a shared pool of packet memory and a dynamic buffer allocation scheme. In this embodiment, the shared memory manager 220 is responsible for managing the shared pool of free buffers in the shared memory 230. It services two categories of clients, buffer consumers (e.g., the input ports 206) and

WO 99/00939

PCT/US98/13365

- 17 -

buffer providers (e.g., the output ports 206). The buffer consumers request free buffers from the shared memory manager 220 at appropriate times during incoming packet reception. Then, during packet forwarding processing, buffer ownership changes hands between the two client types. Finally, at the appropriate times during packet transmission  
5 buffers are returned by the buffer providers to the shared memory manager 220.

Referring now to Figure 3A, a logical view of shared memory 230 is depicted having stored therein packet data in a number of buffers. In this example, the shared memory 230 is segmented into a number of buffers (pages) of programmable size. All the buffers may have the same size, or alternatively, individual buffer sizes may vary. In  
10 another embodiment, the buffers may be further subdivided into a number of memory lines. Each line may be used for storing packet data. In other embodiments, control information may also be associated with each of the memory lines. The control information may include information for facilitating efficient access of the packet data such as an end of packet field. The separation of control information and data increases the  
15 efficiency of accesses to and from the shared memory 230.

A given packet's data may be stored in one or more buffers. In this example, packet #1 is distributed across three buffers 350-352, packet #2's data is stored in three buffers 360-362, and packet #3 is fully contained within one buffer 370. This example also illustrates that the buffers for a particular packet and the packets themselves need not  
20 be in any particular order in the shared memory 230. In this manner, when a particular buffer becomes free, it may be immediately used to fulfill the next buffer request. Also, it may be convenient to limit packet data contained within a particular buffer to one packet. That is, the implementation may be simplified by preventing the mixing of more than one packet within a buffer. In this embodiment, it should be appreciated a packet is

WO 99/00939

PCT/US98/13365

- 18 -

represented as a list of one or more buffers. Therefore, forwarding packet #1 from an input port 206 to an output port 206 would involve removing the pointers to buffers 350-352 from the input port's input queue and transferring them to an output queue of the output port 206.

5

#### EXEMPLARY SHARED MEMORY MANAGER

Figure 3B is a block diagram of the shared memory manager of Figure 2 according to one embodiment of the present invention. According to this embodiment, the shared memory manager 220 includes a buffer tracking unit 329 and a shared memory interface 330. The shared memory interface 330 provides an efficient centralized interface to the shared memory 230. The buffer tracking unit 329 further includes a buffer manager 325. The buffer manager 325 provides a level of indirection which is exploited by the input and output ports 206 by queuing pointers to buffers that contain packet data rather than queuing the packet data itself. As such, the buffering provided by the present invention does not fit into the prior buffering categories such as input packet buffering or output packet buffering. Rather, the buffering architecture described herein is well suited for shared memory buffering with output queuing, for example. Advantageously, since pointers are queued at the ports, the act of forwarding, according to the present embodiment, is simplified to transferring one or more buffer pointers between an input port 206 to an output queue of one or more output ports 206.

Additionally, this flexible approach allows each buffer in the shared memory 230 to be "owned" by one or more different ports at different points in time without having to duplicate the packet data. For example, copies of a multicast packet's buffer pointer(s)

WO 99/00939

PCT/US98/13365

- 19 -

may reside in several output port queues while only one copy of the packet data need reside in the shared memory 230.

The buffer tracking unit 329 additionally includes a pointer random access memory (PRAM) 320. The PRAM 320 may be an on or off-chip pointer table that stores usage counts for buffers of the shared memory 230. Since the assignee of the present invention has found it advantageous to implement each switching element 100 as a single application specific integrated circuit (ASIC), it is preferred that the pointer table be kept compact enough to allow it to be maintained on-chip to facilitate the desired highly integrated implementation.

In any event, with reference to the PRAM 320, the number of buffer owners at a given time is known by the buffer manager 325; thereby allowing the buffer manager 325 to perform efficient real-time free buffer determination for dynamic buffer allocation and allowing efficient deallocation processing of buffers upon their release by the last output port 206. Importantly, the next free buffer, if memory is available, is always kept on hand by the buffer tracking unit 329 for immediate delivery to the requesting input port 206. The processing involved in allocating buffers, transferring buffer ownership, and deallocating buffers will be described in further detail below.

#### EXEMPLARY BUFFER TRACKING PROCESS

Figure 4 is a block diagram of the buffer tracking unit 329 of Figure 3B according to one embodiment of the present invention. In the embodiment depicted, the buffer tracking unit 329 includes an arbiter 470, an array controller 450, an address/data generator 460, PRAM 320, a priority encoder 410, and a pointer generator 440.

WO 99/00939

PCT/US98/13365

- 20 -

According to the present embodiment, the PRAM 320 further includes a count array 430 and a tag array 420. The count array 430 is a memory that stores a count representing the number of ports that are currently using a corresponding buffer in the shared memory 230. In one embodiment, the location of a given count field in the count array 430 represents the start address of the corresponding buffer in the shared memory 230. In this manner, the same pointer may be used to determine the buffer ownership count and to store and retrieve packet data.

In one embodiment, the count array 430 is divided into rows and columns. Each row may store a set of one or more of the plurality of count fields. In this example, the tag array 420 is a memory that has the same number of rows as the count array 430 and contains a field indicating the availability of a buffer in the corresponding row of the count array 430. That is, if any of the count fields in the corresponding row of the count array 430 are zero, meaning no owners, for example, then the tag field is a one, meaning a buffer is available, for example. Advantageously, this indexing mechanism facilitates the real-time indication of free buffers. Alternative configurations are contemplated. For example, in alternative embodiments, the count array 430 and the tag array 420 may share the same memory.

Arbitor 470 arbitrates among the input ports and the output ports 206 to provide only a single port with access to the PRAM 320 at any given time. The arbitor 470 is coupled to the array controller 450 to allow the single port selected to access the PRAM 320. The array controller 450 schedules read and write operations for the PRAM 320 allowing access to both the tag array 420 and the count array 430.

The address/data generator 460 generates control signals for the particular memory or memories employed by the PRAM 320 to facilitate modification of the count fields and

WO 99/00939

PCT/US98/13365

- 21 -

tag fields. Handshake signals for the input and output ports 206 are also generated by the address/data generator 460 as will be described in more detail below. Additionally, the address/data generator 460 may provide a conversion from a buffer pointer to a row address in the count array 430.

5       The priority encoder 410 has inputs corresponding to each element of the tag array 420. In one embodiment, it generates an output which indicates the location of the first non-zero tag bit in the tag array 420. The output of the priority encoder 410 is an input to the pointer generator 440. According to one embodiment, the pointer generator 440 compares the entries from the row indicated by the priority encoder 410 and adds an  
10       encoding representative of the position of an available buffer to produce a buffer pointer for one of the input ports 206.

#### BUFFER ALLOCATION PROCESSING

Figure 5 is a flow diagram illustrating buffer allocation processing according to  
15       one embodiment of the present invention. At step 505, the next free buffer pointer is produced by the pointer generator 440. In one embodiment, the pointer generator 440 attempts to keep one or more pointers available to allow immediate servicing of buffer requests.

At step 510, the count field corresponding to the generated pointer is updated. In  
20       one embodiment, this is accomplished by writing a predetermined value, such as the maximum value, to the count field. For example, the maximum value for a 4-bit counter is 15 or 1111b.

WO 99/00939

PCT/US98/13365

- 22 -

At step 515, if the current row of count fields contains no free buffers after the update of step 510, then at step 520 the tag corresponding to this row is updated to so indicate. Otherwise, processing continues with step 525.

At step 525, the buffer tracking unit 329 waits until one or more input ports 206 request a buffer pointer. Upon detecting one or more requests, processing continues with step 530.

At step 530, one input port request is selected for processing by the buffer tracking unit 329. In one embodiment, the input port requests are received by the arbitor 470. The arbitor 470 selects one of the input port requests for servicing by the buffer tracking unit 329. In another embodiment, the buffer tracking unit 329 may support mixed port speeds by giving priority to the faster network links. For example, the arbitor 470 may be configured to arbitrate between the buffer pointer requests in a prioritized round robin fashion giving priority to the faster interfaces by servicing each slow interface (e.g., Fast Ethernet port) for each N faster interfaces (e.g., Gigabit Ethernet ports).

At step 535, the free buffer pointer is returned to the input port 206 that was selected at step 530. Buffer allocation processing may continue by repeating steps 505-535.

#### BUFFER OWNERSHIP TRANSFER PROCESSING

Figure 6 is a flow diagram illustrating buffer ownership transfer processing according to one embodiment of the present invention. At step 610, an input port 206 determines the number of ports to which a packet is to be forwarded based upon a forwarding decision received from the switch fabric 210.

WO 99/00939

PCT/US98/13365

- 23 -

For each buffer in which the packet's data is stored, the input port 206 performs steps 620-640. At step 620, the input port 206 transfers a buffer pointer to the output port(s) 206 indicated by the forwarding decision. At step 630, the input port 206 notifies the buffer manager 325 of the ownership transfer of the buffer from the input port 206 to the output port(s) 206 by communicating the number of output ports to which the buffer was successfully transferred to the buffer manager 325.

At step 640, the count field associated with the current buffer is updated to reflect the number of output ports that will transmit the buffer. Importantly, the inventors of the present invention have designed the update mechanism described herein to operate in such a manner that does not require the buffer accounting to be race free. Before describing the novel update mechanism, the race condition that is resolved by the update mechanism will now briefly be described.

As should be appreciated, before an input port 206 can notify the buffer manager 325 of the number of output ports to which a particular buffer pointer was transferred, the input port 206 determines if the output port(s) 206 can accept an additional buffer pointer by testing an output queue full indication, for example. It is possible for one or more output ports 206 to receive a buffer pointer, transmit the packet data associated with the buffer pointer, and update the buffer count before the input port 206 has notified the buffer manager 325 of the total number of output ports.

The update mechanism that handles the race condition described above will now be described. According to one embodiment, the buffer manager 325 may be configured to perform a read/modify/write on the count field rather than simply setting the count field to the number indicated by the input port 206. Recall, in the buffer allocation process, according to one embodiment, the count field is set to a predetermined value, such as the



WO 99/00939

PCT/US98/13365

- 24 -

count field's maximum value (e.g., Fh) upon buffer allocation. Therefore, during buffer ownership transfer processing, the count field may be updated to reflect the current number of output ports that will transmit the buffer by reading the current contents of the appropriate count field, adding the number supplied by the input port 206 to the current

5 contents plus a predetermined value to compensate for the initial value written by the buffer tracking unit 329 during buffer pointer allocation, and then writing the result back to the count field. Advantageously, in this manner, the count field will accurately reflect the current number of output ports for the buffer pointer whether or not the count field was previously decremented by one or more output ports 206 as illustrated in Table 1, below.

10 Table 1 illustrates the count field's value after each of the actions in the first column.

- 25 -

Action	Count Field
An input port 206 requests a buffer pointer from the buffer tracking unit 329.	0000b
A buffer pointer is provided to the input port 206.	1111b
The forwarding decision indicates the packet including the buffer is to be forwarded to three output ports 206.	1111b
The input port notifies the buffer tracking unit 329 of the number of owners of the buffer and forwards the buffer pointer to each of the three output ports 206.	1111b
One output port 206 completes transmission of the buffer and notifies the buffer tracking unit 329 that it no longer holds a copy of the buffer pointer.	1111b
The buffer tracking unit 329 processes the output port's notification prior to the input port's notification. Read: 1111b Modify : 1111b - 0001b = 1110b Write: 1110b	1110b
The buffer tracking unit 329 processes the input port's notification which indicates there are 3 buffer owners. Read: 1110b Modify : 1110b + 0011b + 0001b = 0010b Write: 0010b	0010b
The other two output ports 206 complete transmission of the buffer and so notify the buffer tracking unit 329.	0010b
The buffer tracking unit updates the count field.	0000b

Table 1

At step 650, it is determined if all buffers for the packet have been processed. If so, the ownership transfer of this packet is complete; otherwise, processing continues with step 620.

5

#### BUFFER RETURN PROCESSING

Figure 7 is a flow diagram illustrating buffer return processing according to one embodiment of the present invention. After output ports 206 have finished transmitting the contents of particular buffer, the output port 206 returns the buffer pointer so that it may be reused in the buffer allocation processing discussed above.

10

WO 99/00939

PCT/US98/13365

- 26 -

In the present embodiment, at step 710, one or more output ports 206 request to return a buffer. At step 720, the arbitor 470, selects a request to service.

At step 730, the buffer count is updated to reflect the fact that one less output port 206 owns the buffer. For example, the buffer count may be decremented by performing a  
5 read/modify/write operation.

At step 740, if the buffer is now free, processing continues with step 750. A buffer is free when no output ports 206 have a pointer to this buffer pending in any of their output queues. In one embodiment, the buffer is determined to be free based upon the count field being decremented to zero. However, alternative embodiments may use  
10 other indications.

At step 750, a tag corresponding to the set of buffers to which the current buffer belongs is updated to indicate the availability of a buffer in this set of buffers. In one embodiment, a tag array is employed which stores a single bit for each set of buffers.

Having described an exemplary method and apparatus for shared memory  
15 management, the interfaces among the components will now be described.

#### BUFFER MANAGER/INPUT PORT INTERFACE

According to one embodiment, the following signals may be used to implement the handshake between the buffer manager 325 and the input ports 206.

20 (1) Br\_Ptr\_IP - Bus Request for the Input Port Buffer Pointer Data Bus

This signal is asserted by the input ports 206 to the buffer manager 325. At the appropriate point during input packet reception, the input port 206 asserts this signal to indicate to the buffer manager 325 that a buffer pointer is desired. A bus request

WO 99/00939

PCT/US98/13365

- 27 -

acknowledgment (see Br\_Ptr\_IP\_Ack below) is expected to be asserted by the buffer manager 325 in response.

(2) Br\_Ptr\_IP\_Ack - Buffer Pointer Acknowledgment

5 This signal is asserted by the buffer manager 325 to the input port 206 that is to receive the buffer pointer (see Br\_Ptr\_Data\_BM\_to\_IP[X:0] below). This signal is to acknowledge the buffer pointer request (see Br\_Ptr\_IP above). The buffer manager 325 arbitrates between the various requests of the input ports and drives a Bus Request Acknowledgment and the buffer pointer in the same cycle.

10 (3) Br\_Ptr\_Data\_BM\_to\_IP[X:0] - Buffer Manager to Input Port Buffer Pointer Data Bus

This data bus is shared by all input ports 206. It indicates to the input port 206 that received the Bus Request Acknowledgment (see Br\_Ptr\_IP\_Ack above) the buffer pointer to be used for the incoming packet.

(4) Br\_Count - Bus Request for the Count Data Bus

15 This signal is asserted by the input ports 206 to the buffer manager 325. The input port 206 determines the number of output ports that are to receive the packet based upon the forwarding decision received from the switch fabric 210. The input port 206 asserts this signal to indicate to the buffer manager 325 that the number of ports for a buffer pointer is ready. A bus request acknowledgment (see Br\_Count\_Ack below) is expected  
20 to be asserted by the buffer manager 325 in response.

(5) Br\_Count\_Ack - Buffer Count Acknowledgment

This signal is asserted by the buffer manager 325 to the input port 206 that is to provide the number of ports (see Cnt[Y:0] below) for a particular buffer pointer (see Br\_Ptr\_Data\_IP\_to\_BM[X:0] below). This signal is to acknowledge the bus request (see

WO 99/00939

PCT/US98/13365

- 28 -

Br\_Count above) for the count data bus. The buffer manager 325 arbitrates between the various requests of the input ports and drives a bus request acknowledgment to the input port 206 selected by the arbitration.

(6) Dropped\_Ptrs - Number of Ports that Couldn't Receive the Pointers

5 This signal is asserted by the input ports 206 to the buffer manager 325. When the input port 206 cannot post a buffer pointer to all of the output ports 206 indicated by the forwarding decision due to some condition (e.g., full output queue), the input port 206 conveys this information to the buffer manager 325 as it conveys the number of ports. The buffer manager 325 will take this into account when storing the number of output  
10 ports that own the buffer pointer indicated.

(7) Br\_Ptr\_Data\_IP\_to\_BM[X:0] - Input Port to Buffer Manager Buffer Pointer Data Bus

This data bus is shared by all input ports 206. It indicates to the buffer manager 325 the buffer pointer for which the number of ports (see Cnt[Y:0] below) is being  
15 conveyed.

(8) Cnt[Y:0] - Count of Ports

This data bus is shared by all input ports 206. It indicates to the buffer manager 325 the number of ports to which the buffer pointer (see Br\_Ptr\_Data\_IP\_to\_BM[X:0] above) has been transferred.  
20

#### BUFFER MANAGER/OUTPUT PORT INTERFACE

According to one embodiment, the following signals may be used to implement the handshake between the buffer manager 325 and the output ports 206.

(1) Br\_Ptr\_OP - Bus Request for the Output Port Buffer Pointer Data Bus

- 29 -

This signal is asserted by the output ports 206 to the buffer manager 325. At the appropriate point during output packet processing, the output port 206 asserts this signal to indicate to the buffer manager 325 that a buffer pointer is being returned. A bus request acknowledgment (see Br\_Ptr\_OP\_Ack below) is expected to be asserted by the buffer manager 325 in response.

(2) Br\_Ptr\_Data\_OP\_to\_BM[X:0] - Output Port to Buffer Manager Buffer Pointer Data Bus

This data bus is shared by all output ports 206. It indicates to the buffer manager 325 the buffer pointer that is being returned. Output ports 206 return buffer pointers after the data stored in the corresponding buffer has been transmitted.

(3) Br\_Ptr\_OP\_Ack - Buffer Request Acknowledgment

This signal is asserted by the buffer manager 325 to the output port 206 that is to return their buffer pointer (see Br\_Ptr\_Data\_OP\_to\_BM[X:0] above). This signal is to acknowledge the bus request (see Br\_Ptr\_OP above). The buffer manager 325 arbitrates between the various requests of the output ports 206 and drives a Bus Request Acknowledgment to the output port 206 selected by the arbitration logic.

#### INPUT PORT/OUTPUT PORT INTERFACE

According to one embodiment, the following signals may be used to transfer packet ownership from the input ports 206 to the output ports 206.

(1) Arb\_OP\_Ptr - Arbitrated Output Port Buffer Pointer Data Bus

This multiplexed data bus is driven by the output bus arbiter. It is shared by all the output ports 206 for the transfer of buffer pointer ownership information.

(2) OP\_Queue\_Full - Output Port Queue Full

WO 99/00939

PCT/US98/13365

- 30 -

This signal is asserted by the output ports 206 to the input ports 206. This signal is used by the input ports 206 to make filtering decisions when broadcasting packet pointers. That is, if the forwarding decision indicates a packet is to be forwarded to a given output port 206 and that output port's queue is full, then the packet pointer will not  
5 be transferred to that output port 206 and the buffer manager 325 may be notified of the dropped packet pointer (see Dropped\_Ptrs above). Alternatively, the buffer manager 325 may simply be notified of the total number of output ports that have been provided with a particular packet pointer.

For the sake of example, only one output queue has been assumed, note however,  
10 in alternative embodiments more than one output queue may be employed for each output port 206. In this case, a queue full indication may be provided for each additional output queue.

Thus, a buffered architecture has been described which provides temporary storage of received packets in a shared pool of packet memory and provides for efficient allocation  
15 of per port buffering that is proportional to the amount of traffic through a given port.

In the foregoing specification, the invention has been described with reference to specific embodiments thereof. It will, however, be evident that various modifications and changes may be made thereto without departing from the broader spirit and scope of the  
20 invention. The specification and drawings are, accordingly, to be regarded in an illustrative rather than a restrictive sense.

---

- 31 -

CLAIMS

What is claimed is:

- 1 1. A shared memory manager for use in a packet forwarding device, the shared  
2 memory manager comprising:  
3 a pointer memory configured to store information regarding buffer usage in a  
4 shared memory;  
5 an encoder coupled to the pointer memory, the encoder configured to generate an  
6 output which indicates a set of buffers of a plurality of sets of buffers that  
7 contains one or more free buffers; and  
8 a pointer generator coupled to the encoder, the pointer generator locating a free  
9 buffer in the set of buffers and producing a pointer to the free buffer based  
10 upon the output of the encoder and the free buffer's location within the set  
11 of buffers.
- 1 2. The shared memory manager of claim 1, wherein the information regarding buffer  
2 usage comprises usage counts for each of a plurality of buffers.
- 1 3. The shared memory manager of claim 2, wherein the pointer memory further  
2 comprises a count array including a plurality of entries configured to store the  
3 usage counts, each entry of the plurality of entries corresponding to one of the  
4 plurality of buffers, the usage counts representing the number of ports that hold a  
5 pointer to the corresponding buffer.



WO 99/00939

PCT/US98/13365

- 32 -

- 1 4. The shared memory manager of claim 3, wherein a given entry's location in the  
2 count array represents the address of the corresponding buffer in the shared  
3 memory.
- 1 5. The shared memory manager of claim 3, wherein the pointer memory further  
2 comprises a tag array coupled to the count array, the tag array including an  
3 indication corresponding to each set of buffers in the plurality of sets of buffers,  
4 the indication indicating the availability of one or more buffers in the  
5 corresponding set of buffers.
- 1 6. A packet forwarding device comprising:  
2 a plurality of output ports configured to transmit packets onto a network;  
3 a plurality of input ports coupled to the plurality of output ports, the plurality of  
4 input ports configured to receive a packet from the network and forward  
5 the packet to one or more of the plurality of output ports;  
6 a shared memory coupled to the plurality of output ports and the plurality of input  
7 ports, the shared memory segmented into a plurality of buffers configured  
8 to buffer the packet and further configured to store no more than one copy  
9 of a given packet in the shared memory at any given time; and  
10 a shared memory manager coupled to the plurality of input ports and to the  
11 plurality of output ports, the shared memory manager dynamically  
12 allocating buffers from the plurality of buffers to the plurality of input ports  
13 based upon requests received from the plurality of input ports, the shared  
14 memory manager tracking ownership counts for each of the plurality of

WO 99/00939

PCT/US98/13365

- 33 -

15 buffers based upon information received from the plurality of input ports  
16 and the plurality of output ports.

1 7. The packet forwarding device of claim 6, wherein the plurality of input ports are  
2 configured to forward a packet to one or more of the plurality of output ports by  
3 transferring one or more pointers to the packet to one or more of the plurality of  
4 output ports.

1 8. A method of packet forwarding comprising the steps of:  
2 dynamically allocating one or more buffers in a shared memory by determining one  
3 or more free buffer pointers, each of the one or more free buffer pointers  
4 corresponding to one of the one or more buffers, wherein the allocation is  
5 unconstrained by of the locations of the one or more buffers within the  
6 shared memory;  
7 receiving a packet from a first attached network segment;  
8 storing the packet in the one or more buffers;  
9 transferring ownership of the one or more buffer pointers from an input port to one  
10 or more output ports based upon a forwarding decision;  
11 retrieving the packet from the one or more buffers; and  
12 transmitting the packet onto a second attached network segment.

1 9. The method of claim 8, wherein the step of dynamically allocating one or more  
2 buffers in a shared memory by determining one or more free buffer pointers  
3 further includes the step of updating a usage count corresponding to each of the  
4 one or more free buffer pointers.

WO 99/00939

PCT/US98/13365

- 34 -

1 10. The method of claim 9, wherein the step of updating a usage count corresponding  
2 to the free buffer pointer comprises the step of setting the usage count to a  
3 predetermined value to accommodate a potential race condition in usage count  
4 processing.

1 11. The method of claim 8, wherein the step of transferring ownership of the one or  
2 more buffer pointers from an input port to one or more output ports based upon a  
3 forwarding decision further includes the steps of:  
4 for each buffer of the one or more buffers  
5 performing a dequeue operation to remove the corresponding buffer pointer  
6 from an input queue;  
7 performing an enqueue operation to insert the buffer pointer into an output  
8 queue for one or more output ports indicated by the forwarding  
9 decision,  
10 notifying the shared memory manager of the number of output ports to  
11 which the buffer pointer has been successfully enqueued, and  
12 updating a usage count corresponding to the buffer pointer.

1 12. The method of claim 11, wherein the step of updating a usage count corresponding  
2 to the buffer pointer comprises the steps of:  
3 determining a current value of the usage count;  
4 modifying the current value in a manner that accounts for buffers that may have  
5 been deallocated before the step of notifying the shared memory manager  
6 of the number of output ports; and

WO 99/00939

PCT/US98/13365

- 35 -

7 replacing the usage count with the modified value, the modified value reflecting the  
8 number of output ports that currently hold a copy of the buffer pointer;  
9 whereby the adverse effects of race conditions are avoided by accounting for  
10 buffers that may have been deallocated in the modifying step.

1 13. The method of claim 11, wherein the forwarding decision identifies a set of one or  
2 more output ports, and wherein the method further includes the step of determining  
3 a subset of the set of one or more output ports to which to transfer the one or more  
4 buffer pointers.

1 14. The method of claim 13, further including the step of the one or more output ports  
2 generating queue status indications, and wherein the step of determining a subset  
3 of the set of one or more output ports to which to transfer the one or more buffer  
4 pointers is based upon the queue status indications generated by the one or more  
5 output ports.

1 15. The method of claim 8, wherein the step of transmitting the packet onto a second  
2 attached network segment further includes the steps of:  
3 for each buffer of the one or more buffers  
4 transmitting packet data from the buffer onto the second attached network  
5 segment, and  
6 deallocating the corresponding buffer pointer by returning the buffer  
7 pointer to a centralized buffer manager, whereby the buffer  
8 becomes available for storing packet data from another received  
9 packet.

WO 99/00939

PCT/US98/13365

- 36 -

- 1 16. The method of claim 8, wherein no more than one copy of the packet is stored in  
2 the shared memory.
- 1 17. A method of port mirroring comprising the steps of:  
2 receiving a packet;  
3 storing packet data from the packet in one or more buffers;  
4 receiving a forwarding decision, the forwarding decision indicating an output port  
5 and one or more mirror ports; and  
6 transferring the one or more buffer pointers to the output port and to the one or  
7 more mirror ports, whereby packets transmitted on the output port are  
8 mirrored on the one or more mirror ports.
- 1 18. A machine-readable medium having stored thereon data representing sequences of  
2 instructions, said sequences of instructions which, when executed by a processor,  
3 cause said processor to perform the steps of:  
4 dynamically allocating one or more buffers in a shared memory by determining one  
5 or more free buffer pointers, each of the one or more free buffer pointers  
6 corresponding to one of the one or more buffers, wherein the allocation is  
7 unconstrained by of the locations of the one or more buffers within the  
8 shared memory;  
9 receiving a packet from a first attached network segment;  
10 storing the packet in the one or more buffers;  
11 transferring ownership of the one or more buffer pointers from an input port to one  
12 or more output ports based upon a forwarding decision;  
13 retrieving the packet from the one or more buffers; and

WO 99/00939

PCT/US98/13365

- 37 -

14 transmitting the packet onto a second attached network segment.

1 19. The machine-readable medium of claim 18, wherein the step of dynamically  
2 allocating one or more buffers in a shared memory by determining one or more  
3 free buffer pointers the step of updating a usage count corresponding to each of the  
4 one or more free buffer pointers.

1 20. The machine-readable medium of claim 18, wherein the step of transferring  
2 ownership of the one or more buffer pointers from an input port to one or more  
3 output ports based upon a forwarding decision further includes the steps of:  
4 for each buffer of the one or more buffers  
5 performing a dequeue operation to remove the corresponding buffer pointer  
6 from an input queue;  
7 performing an enqueue operation to insert the buffer pointer into an output  
8 queue for one or more output ports indicated by the forwarding  
9 decision,  
10 notifying the shared memory manager of the number of output ports to  
11 which the buffer pointer has been successfully enqueued, and  
12 updating a usage count corresponding to the buffer pointer.

1 21. The machine-readable medium of claim 18, wherein the step of transmitting the  
2 packet onto a second attached network segment further includes the steps of:  
3 for each buffer of the one or more buffers  
4 transmitting packet data from the buffer onto the second attached network  
5 segment, and

WO 99/00939

PCT/US98/13365

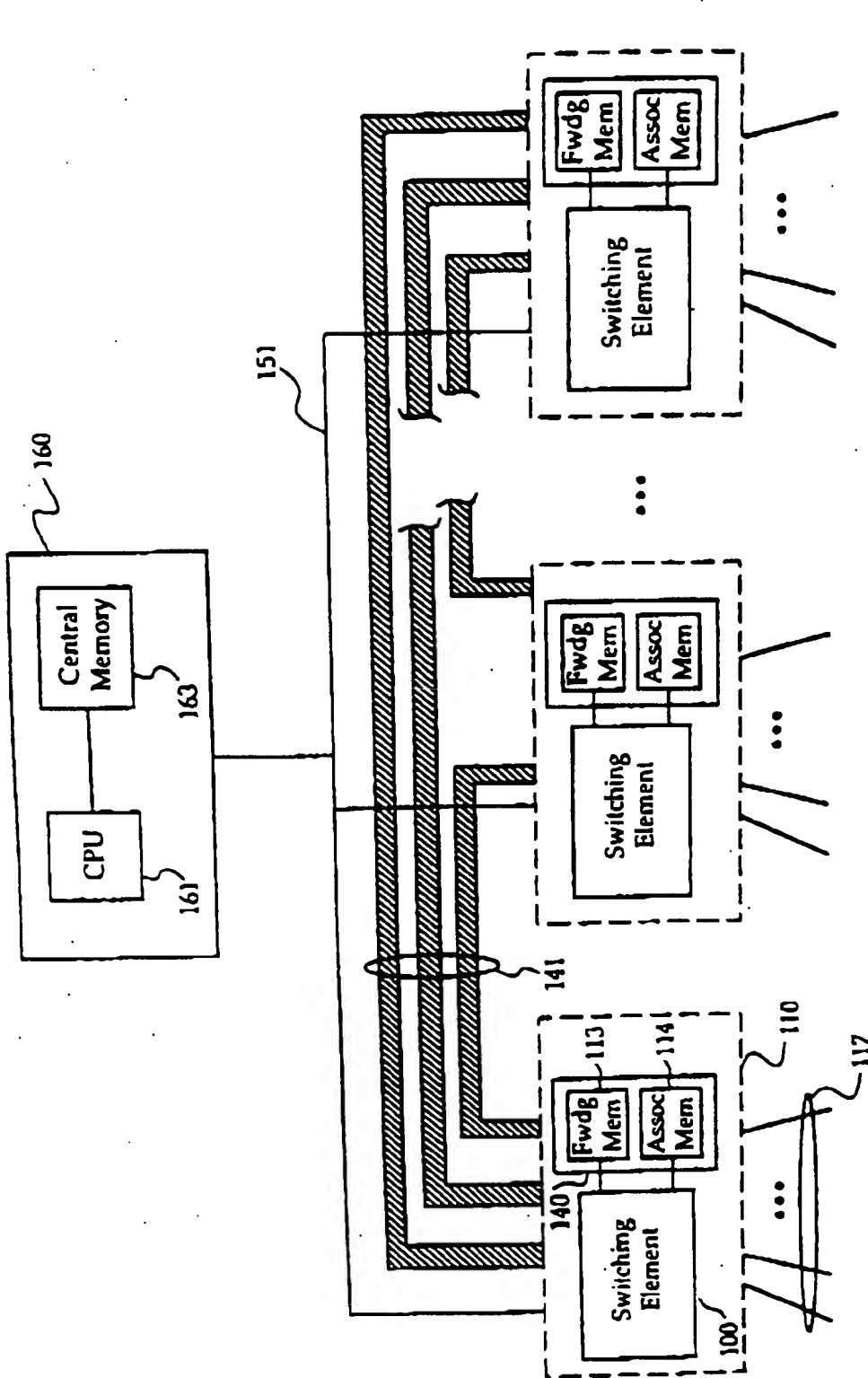
- 38 -

6           deallocating the corresponding buffer pointer by returning the buffer  
7           pointer to a centralized buffer manager, whereby the buffer  
8           becomes available for storing packet data from another received  
9           packet.

WO 99/00939

PCT/US98/13365

1 / 7



To nodes and end-stations

FIG. 1

101



WO 99/00939

PCT/US98/13365

2 / 7

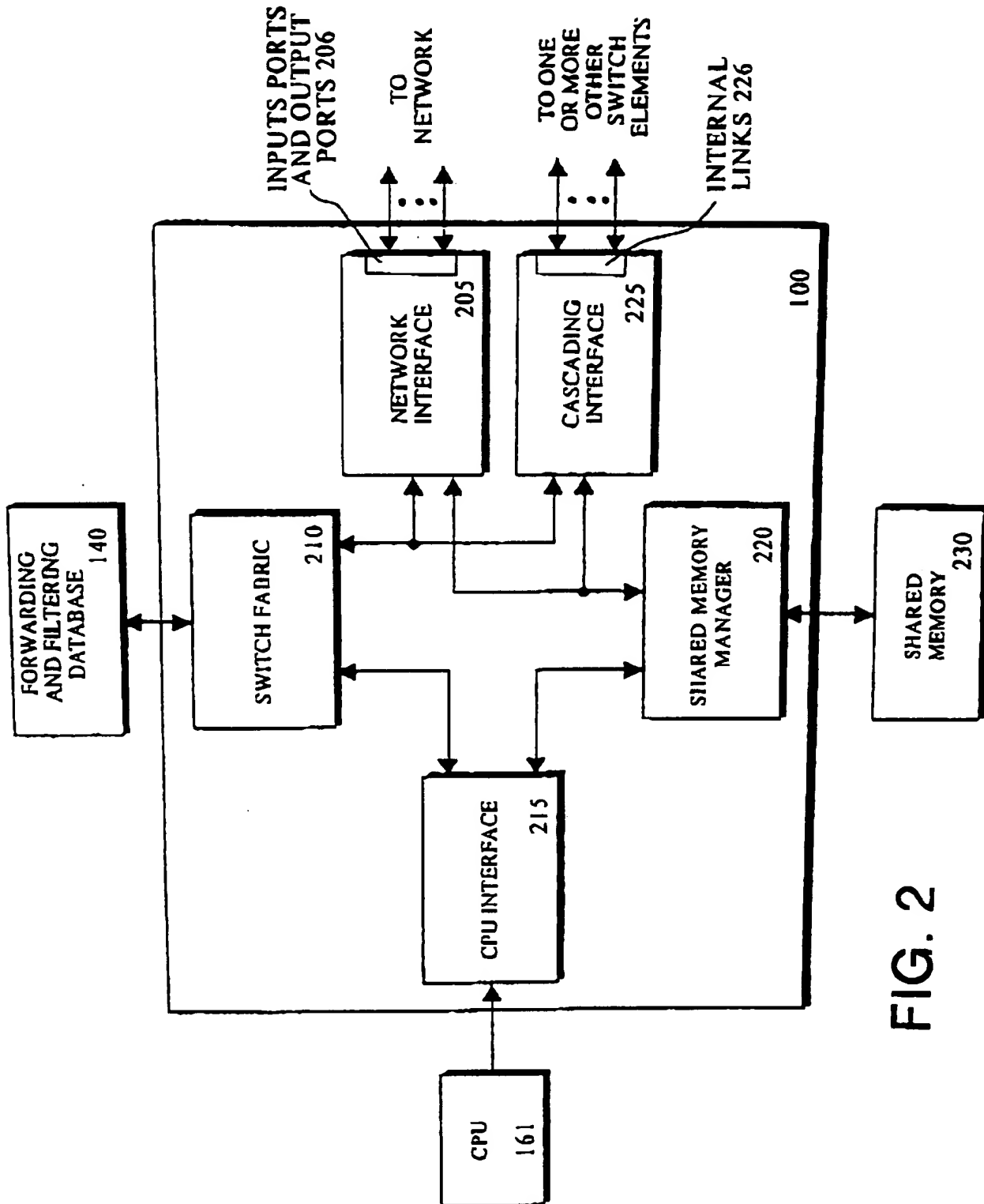


FIG. 2

WO 99/00939

PCT/US98/13365

3/7

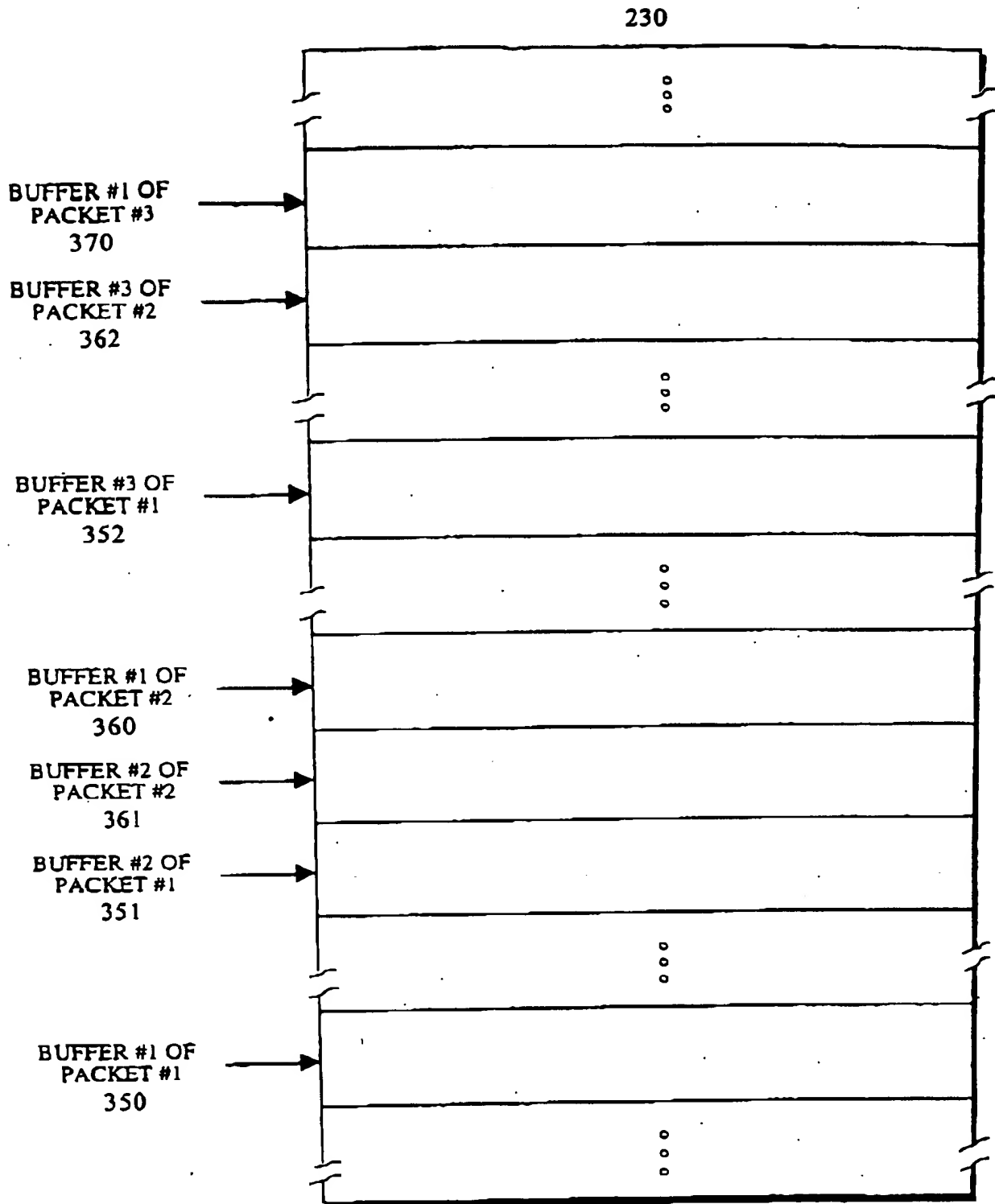


FIG. 3A

WO 99/00939

PCT/US98/13365

4 / 7

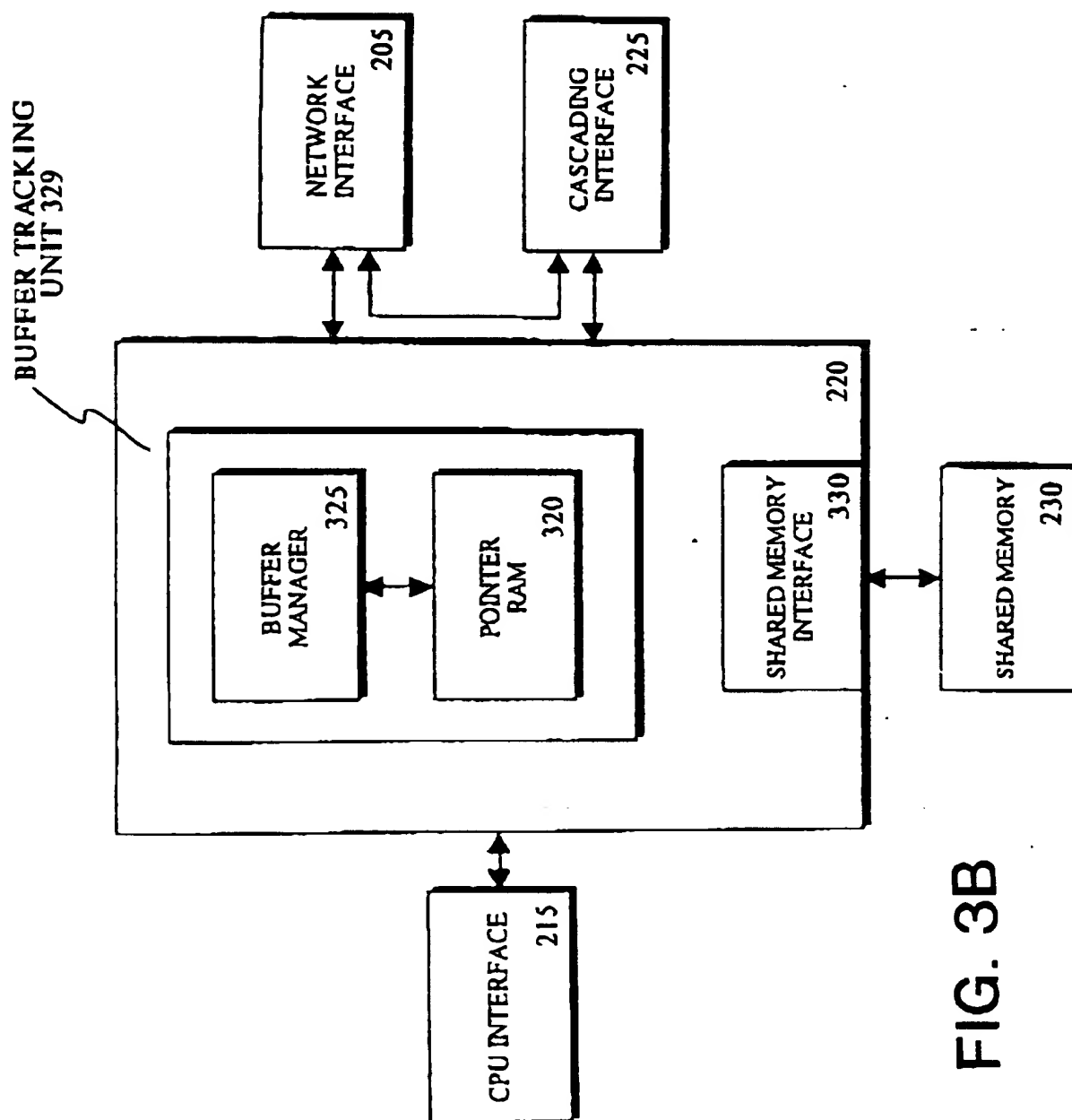


FIG. 3B

5 / 7

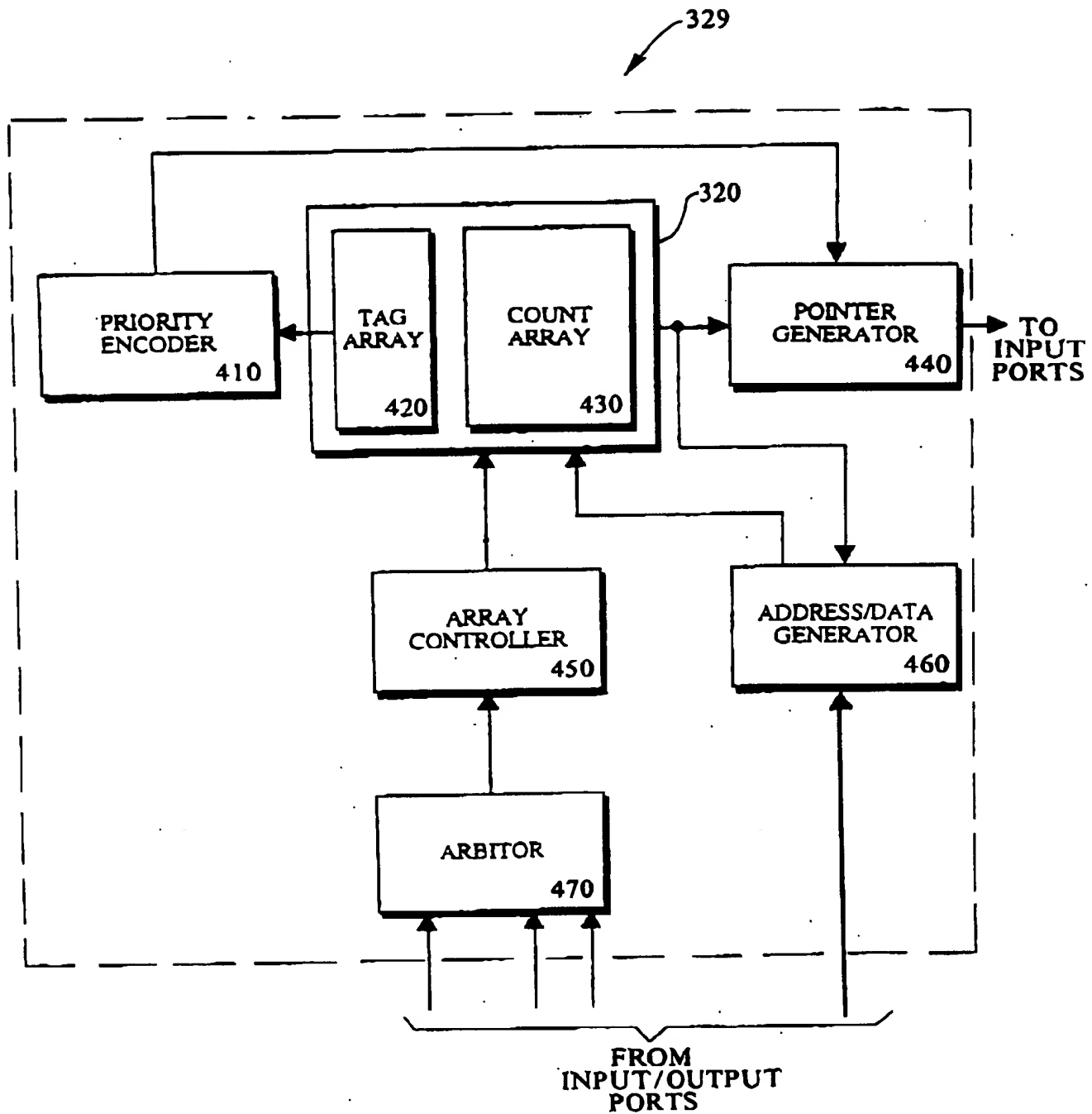


FIG. 4

WO 99/00939

PCT/US98/13365

6 / 7

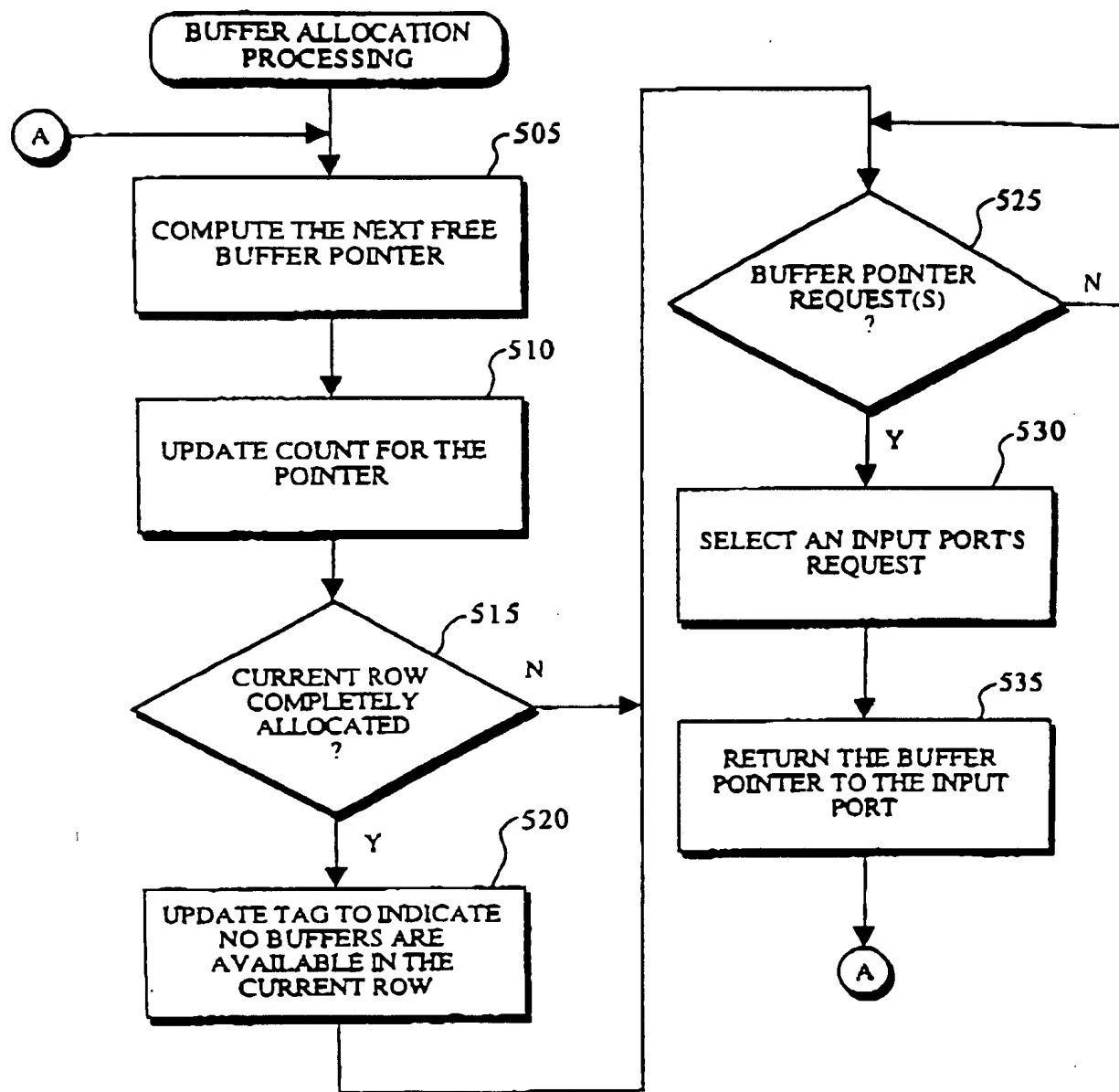
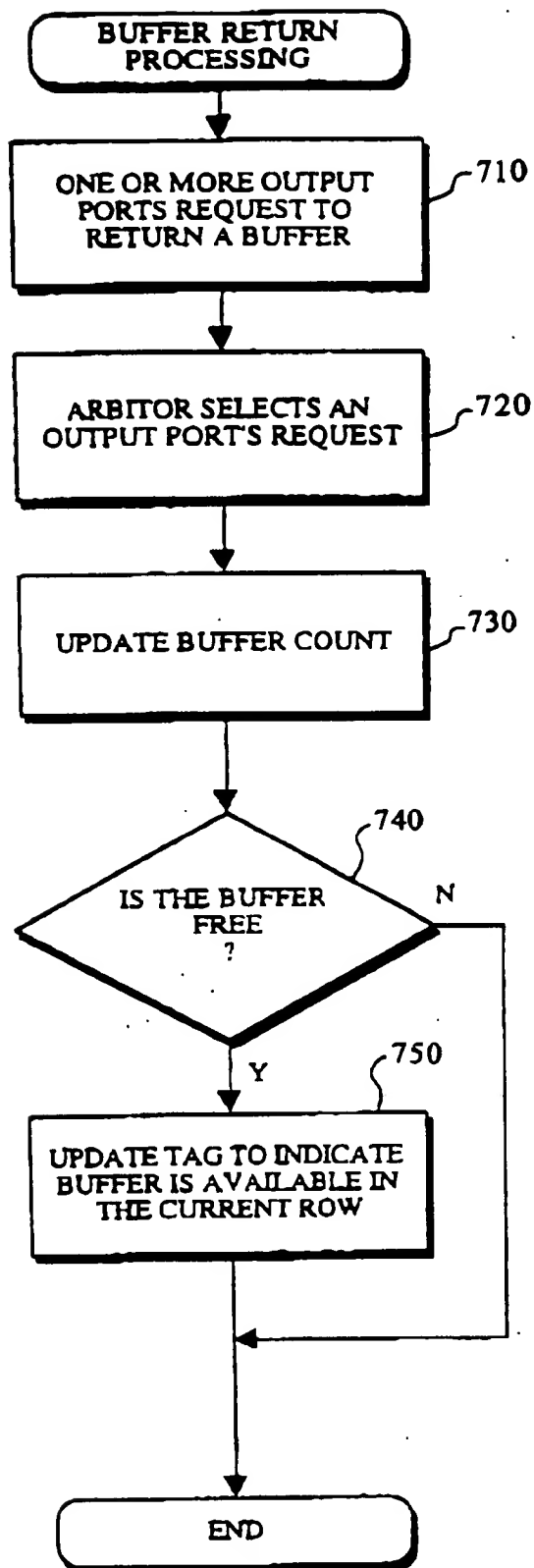
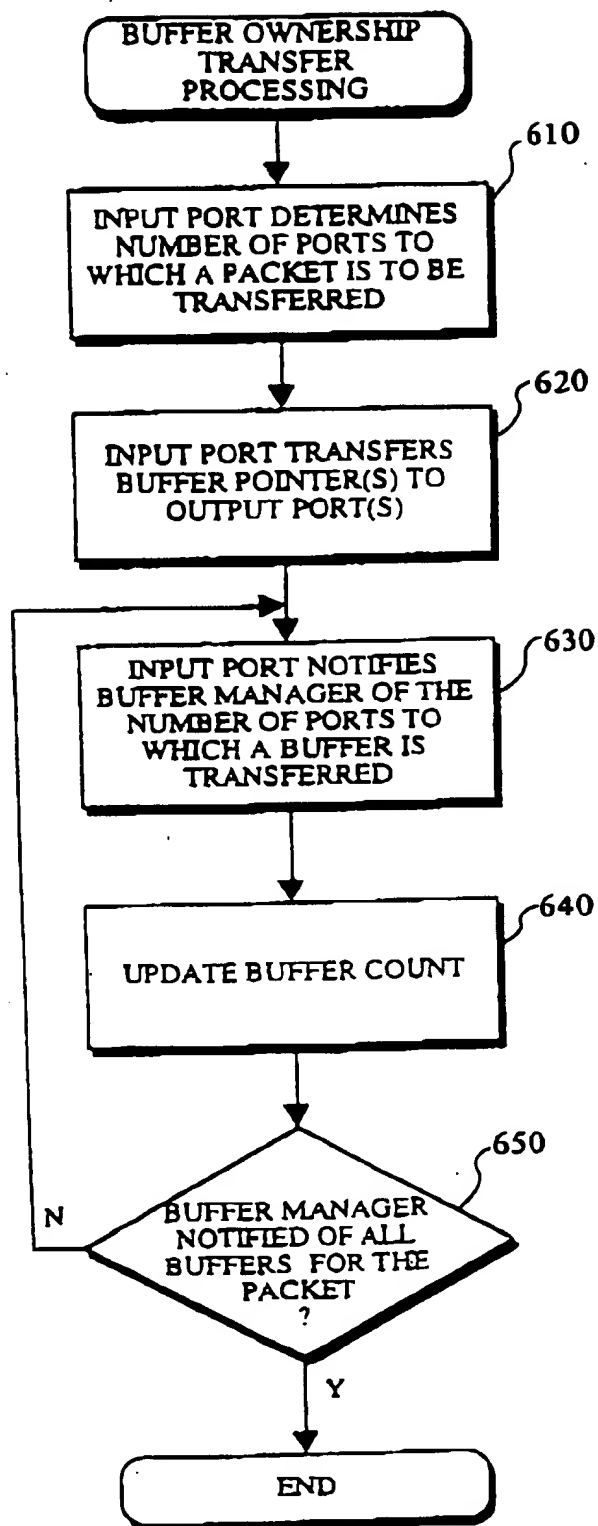


FIG. 5

WO 99/00939

PCT/US98/13365

7 / 7



## INTERNATIONAL SEARCH REPORT

International application No.  
PCT/US98/13365

## A. CLASSIFICATION OF SUBJECT MATTER

IPC(6) : H04L 12/28, 12/56

US CL : 370/351, 390, 395, 411, 412

According to International Patent Classification (IPC) or to both national classification and IPC

## B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

U.S. : 370/351, 390, 395, 411, 412

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)

## C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X --- Y	US 5,633,865 A (Short) 27 May 1997 (27.05.97), col. 6, line 16 to col. 7, line 46.	1-5 --- 6-21
Y	US 5,602,841 A (Lebizay et al.) 11 Feb 1997 (11.02.97) see col. 23, line 51 to col. 24, line 4.	6-7
Y	US 5,651,002 A (Van Seters et al.) 22 Jul 1997 (22.07.97), see FIG. 4B, col. 9, line 18 to col. 11, line 39.	8-21



Further documents are listed in the continuation of Box C.



See patent family annex.

* Special categories of cited documents:	*T	later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention
*A* document defining the general state of the art which is not considered to be of particular relevance	*X*	document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone
*E* earlier document published on or after the international filing date	*Y*	document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art
*L* document which may throw doubt on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)	*Z*	document member of the same patent family
*O* document referring to an oral disclosure, use, exhibition or other means		
*P* document published prior to the international filing date but later than the priority date claimed		

Date of the actual completion of the international search

18 SEPTEMBER 1998

Date of mailing of the international search report

23 OCT 1998

Name and mailing address of the ISA/US  
Commissioner of Patents and Trademarks  
Box PCT  
Washington, D.C. 20231

Facsimile No. (703) 305-3230

Authorized officer

DUONG, FRANK

Telephone No. (703) 308-5428